

## Choice Models in Operations

### Lecture 11: EM algorithm and Frank-Wolfe algorithm

Instructor: Srikanth Jagabathula

Scribe: JinHyun Kim

Last week: Derived the MM algorithm for estimating the parameters of the MNL model through maximizing the log-likelihood function.

Necessary & Sufficient condition for the existence of a unique & bounded optimal solution to

$$\max_{\underline{\beta}} \sum_{t=1}^T \log \frac{e^{\underline{\beta}^T \underline{z}_{j_t,t}}}{\sum_{i \in S_t} e^{\underline{\beta}^T \underline{z}_{i,t}}}$$

is that  $\underline{\gamma}^T(\underline{z}_{i,t} - \underline{z}_{j_t,t}) \leq 0 \quad \forall i \in S_t, j_t, t = 1, 2, \dots, T \implies \underline{\gamma} = \underline{0}$ .

proof: [Necessity] Suppose  $\exists$  a  $\underline{\gamma} \neq \underline{0}$  such that  $\underline{\gamma}^T(\underline{z}_{i,t} - \underline{z}_{j_t,t}) \leq 0 \quad \forall i \in S_t, j_t, t$ . Consider

$$l_t(\underline{\beta}) = \log \frac{e^{\underline{\beta}^T \underline{z}_{j_t,t}}}{\sum_{i \in S_t} e^{\underline{\beta}^T \underline{z}_{i,t}}} = -\log \sum_{i \in S_t} e^{\underline{\beta}^T(\underline{z}_{i,t} - \underline{z}_{j_t,t})}$$

Then, we have for any  $c > 0$ ,

$$l_t(\underline{\beta} - c\underline{\gamma}) = \log \sum_{i \in S_t} e^{\underline{\beta}^T(\underline{z}_{i,t} - \underline{z}_{j_t,t}) - c\underline{\gamma}^T(\underline{z}_{i,t} - \underline{z}_{j_t,t})} = -\log(1 + \sum_{i \in S_t \setminus \{j_t\}} e^{\underline{\beta}^T(\underline{z}_{i,t} - \underline{z}_{j_t,t}) - c\underline{\gamma}^T(\underline{z}_{i,t} - \underline{z}_{j_t,t})})$$

We have 2 cases:

(i)  $\exists$  an  $i, j_t$  such that  $\underline{\gamma}^T(\underline{z}_{i,t} - \underline{z}_{j_t,t}) < 0$ .

Let  $g_t(c) = l_t(\underline{\beta} - c\underline{\gamma})$ .

$$g'_t(c) = \frac{\sum_{i \in S_t \setminus \{j_t\}} \underline{\gamma}^T(\underline{z}_{i,t} - \underline{z}_{j_t,t}) e^{\underline{\beta}^T(\underline{z}_{i,t} - \underline{z}_{j_t,t}) - c\underline{\gamma}^T(\underline{z}_{i,t} - \underline{z}_{j_t,t})}}{1 + \sum_{i \in S_t \setminus \{j_t\}} e^{\underline{\beta}^T(\underline{z}_{i,t} - \underline{z}_{j_t,t}) - c\underline{\gamma}^T(\underline{z}_{i,t} - \underline{z}_{j_t,t})}}$$

Because the derivative w.r.t  $c$  is always  $< 0$ , decreasing the value of  $c$  will strictly increase the value of  $l_t(\underline{\beta} - c\underline{\gamma}) \implies$  the optimal solution cannot be bounded.

(ii) Suppose  $\underline{\gamma}^T(\underline{z}_{i,t} - \underline{z}_{j_t,t}) = 0 \quad \forall i \in S_t, j_t, t = 1, 2, \dots, T$ .

$$\begin{aligned} \implies l(\underline{\beta} - c\underline{\gamma}) &= \sum_t l_t(\underline{\beta} - c\underline{\gamma}) = \sum_t l_t(\underline{\beta}) = l(\underline{\beta}) \quad \forall c \in \mathbb{R} \\ \implies &\text{multiple optima} \end{aligned}$$

## EM algorithm as a special case of MM algorithm

EM setup: Suppose we observe a data point  $\underline{y}$  according to some distribution parametrized by  $\underline{\theta}$ .

Goal: Estimate  $\underline{\theta}$  via MLE, i.e. solve  $\max_{\underline{\theta}} \log p(\underline{y}|\underline{\theta})$  where  $p(\cdot|\underline{\theta})$  is the distribution according to which  $\underline{y}$  was generated.

Suppose there is a latent variable such that if the variable is observed, then the MLE problem becomes simple(r). Let  $\underline{z}$  denote the latent variable and let  $p(\underline{y}, \underline{z}|\underline{\theta})$  denote the joint probability distribution function.

Then, the complete-data log-likelihood function is defined as

$$l_c(\underline{y}, \underline{z}) = \log p(\underline{y}, \underline{z}|\underline{\theta})$$

Correspondingly, we refer to  $\log p(\underline{y}|\underline{\theta})$  as the incomplete-data log-likelihood function, denoted by  $l_{IC}(\underline{y})$ .

$$\max_{\underline{\theta}} l_{IC}(\underline{y}; \underline{\theta}) = \max_{\underline{\theta}} \log p(\underline{y}|\underline{\theta}) = \max_{\underline{\theta}} \log \sum_{\underline{z}} p(\underline{y}, \underline{z}|\underline{\theta})$$

Let's apply the MM meta algorithm.

Let  $\underline{\theta}^{(t)}$  be the current iterate.

The, we want to find  $g(\cdot|\underline{\theta}^{(t)})$  that is a minorizing function, i.e.  $l_{IC}(\underline{y}; \underline{\theta}) \geq g(\underline{\theta}|\underline{\theta}^{(t)})$  and  $l_{IC}(\underline{y}; \underline{\theta}^{(t)}) = g(\underline{\theta}^{(t)}|\underline{\theta}^{(t)})$ .

Trick: Suppose we have a function  $f(\cdot)$  that is strict concave and say our goal is to find a minorizing function for  $f(\sum_i x_i)$  at the current iterate  $\underline{x}^{(t)}$ . Since  $f(\cdot)$  is concave, we must have  $f(\sum_i \alpha_i y_i) \leq \sum_i \alpha_i f(y_i) \quad \forall \alpha_i \geq 0 \quad \forall i, \sum \alpha_i = 1$ . Equality holds iff  $y_i = y_j \quad \forall i, j$ .

Set  $\alpha_i = \frac{x_i^{(t)}}{\sum_j x_j^{(t)}}$  and  $y_i = \frac{x_i}{x_i^{(t)}} \sum_j x_j^{(t)}$ .

Suppose  $x_i^{(t)} > 0 \quad \forall i$  and  $\sum_j x_j^{(t)} > 0$

$$\implies f(\sum_i x_i) = f(\sum_i \alpha_i y_i) \geq \sum_i \alpha_i f(y_i) = \sum_i \frac{x_i^{(t)}}{\sum_j x_j^{(t)}} f(x_i \frac{\sum_j x_j^{(t)}}{x_i^{(t)}})$$

equality occuring iff  $\frac{x_i}{x_i^{(t)}} = \frac{x_j}{x_j^{(t)}} \quad \forall i, t$ .

Therefore,  $g(\underline{x}|\underline{x}^{(t)}) \triangleq \sum_i x_i^{(t)} f(x_i \frac{\sum_j x_j^{(t)}}{x_i^{(t)}})$  is a minorizing function to  $f(\sum_i x_i)$  at  $\underline{x}^{(t)}$ .

Now, we apply this general idea to our setting by taking  $f(\cdot)$  to be  $\log(\cdot)$  and  $x_{\underline{z}} = p(\underline{y}, \underline{z}|\underline{\theta}) \quad \forall \underline{z}$ .

$$\begin{aligned} \therefore \log(\sum_{\underline{z}} p(\underline{y}, \underline{z}|\underline{\theta})) &\geq \sum_{\underline{z}} \frac{p(\underline{z}, \underline{y}|\underline{\theta}^{(t)})}{\sum_{\underline{z}'} p(\underline{z}', \underline{y}|\underline{\theta}^{(t)})} \log[p(\underline{y}, \underline{z}|\underline{\theta}) \frac{\sum_{\underline{z}'} p(\underline{z}', \underline{y}|\underline{\theta}^{(t)})}{p(\underline{z}, \underline{y}|\underline{\theta}^{(t)})}] \\ &= \sum_{\underline{z}} \frac{p(\underline{z}, \underline{y}|\underline{\theta}^{(t)})}{\sum_{\underline{z}'} p(\underline{z}', \underline{y}|\underline{\theta}^{(t)})} \log p(\underline{y}, \underline{z}|\underline{\theta}) + \text{constant} = \sum_{\underline{z}} \frac{p(\underline{z}, \underline{y}|\underline{\theta}^{(t)})}{p(\underline{y}|\underline{\theta}^{(t)})} \log p(\underline{z}, \underline{y}|\underline{\theta}) + \text{constant} \\ &= \sum_{\underline{z}} p(\underline{z}|\underline{y}, \underline{\theta}^{(t)}) \log p(\underline{z}, \underline{y}|\underline{\theta}) + \text{constant} = \mathbb{E}_{\underline{Z}|\underline{y}, \underline{\theta}^{(t)}} [\log p(\underline{Z}, \underline{y}|\underline{\theta})] \\ &\longrightarrow \text{E-step (expectation-step)} \end{aligned}$$

We then maximize the expectation above:

$$\underline{\theta}^{(t+1)} = \underset{\underline{\theta}}{\operatorname{argmax}} \mathbb{E}_{\underline{Z}|\underline{y}, \underline{\theta}^{(t)}} [\log p(\underline{Z}, \underline{y}|\underline{\theta})]$$

→ M-step (Max-step)

We use the EM framework to derive the estimation algorithm for a latent-class MNL model with K classes.

Data: We observe choice from m customers. For customer  $i$ , we observe choices  $(j_{i,t}, S_{i,t})$  for  $t \in T_i$ . The log-likelihood function is

$$\underline{\alpha}, \underline{\beta}: \sum_{\alpha_k=1, \alpha_k \geq 0 \forall k} \max \sum_i \log \left[ \sum_{k=1}^K \alpha_k \prod_{t \in T_i} \frac{e^{\beta_{k,j_{i,t}}}}{\sum_{l \in S_{i,t}} e^{\beta_{k,l}}} \right]$$

We introduce latent variable  $z_i$  which is the class membership of each customer  $i$ . The complete-data log-likelihood function can be written as

$$\begin{aligned} l_c(z, \underline{y}; \underline{\theta}) &= \sum_i \sum_k \mathbb{1}[z_i = k] \log \left[ \left( \prod_{t \in T_i} \frac{e^{\beta_{k,j_{i,t}}}}{\sum_{l \in S_{i,t}} e^{\beta_{k,l}}} \right) \alpha_k \right] \\ &= \sum_i \sum_k \mathbb{1}[z_i = k] \left[ \log \alpha_k + \sum_{t \in T_i} \log \frac{e^{\beta_{k,j_{i,t}}}}{\sum_{l \in S_{i,t}} e^{\beta_{k,l}}} \right] \end{aligned}$$

Current iterate  $\underline{\alpha}^{(t)}, \underline{\beta}^{(t)}$ .

E-step:

$$\begin{aligned} \mathbb{E}_{\underline{z}|\underline{y}, \underline{\theta}^{(t)}} [l_c(z, \underline{y}; \underline{\theta})] &= \mathbb{E}_{\underline{z}|Data, \underline{\alpha}^{(t)}, \underline{\beta}^{(t)}} \left[ \sum_i \sum_k \mathbb{1}[z_i = k] \left[ \log \alpha_k + \sum_{t \in T_i} \log \frac{e^{\beta_{k,j_{i,t}}}}{\sum_{l \in S_{i,t}} e^{\beta_{k,l}}} \right] \right] \\ &= \sum_i \sum_{k=1}^K h_{ik}^{(t)} \left[ \log \alpha_k + \sum_{t \in T_i} \log \frac{e^{\beta_{k,j_{i,t}}}}{\sum_{l \in S_{i,t}} e^{\beta_{k,l}}} \right] \end{aligned}$$

where  $h_{ik}^{(t)} = \mathbb{E}_{\underline{z}|Data, \underline{\alpha}^{(t)}, \underline{\beta}^{(t)}} [\mathbb{1}[z_i = k]] = Pr(z_i = k | Data, \underline{\alpha}^{(t)}, \underline{\beta}^{(t)})$

$$\begin{aligned} &= \frac{Pr(Data_i | z_i = k, \underline{\alpha}^{(t)}, \underline{\beta}^{(t)}) Pr(z_i = k)}{\sum_{k'} Pr(Data_i | z_i = k', \underline{\alpha}^{(t)}, \underline{\beta}^{(t)}) Pr(z_i = k')} \\ &= \frac{\alpha_k^{(t)} \prod_{t \in T_i} (e^{\beta_{k,j_{i,t}}^{(t)}} / \sum_{l \in S_{i,t}} e^{\beta_{k,l}^{(t)}})}{\sum_{k'} \alpha_{k'}^{(t)} \prod_{t \in T_i} (e^{\beta_{k',j_{i,t}}^{(t)}} / \sum_{l \in S_{i,t}} e^{\beta_{k',l}^{(t)}})} \end{aligned}$$

M-step:

$$\underline{\alpha}, \underline{\beta}: \sum_{\alpha_k=1, \alpha_k \geq 0 \forall k} \max \sum_i \sum_{k=1}^K h_{ik}^{(t)} \left[ \log \alpha_k + \sum_{t \in T_i} \log \frac{e^{\beta_{k,j_{i,t}}}}{\sum_{l \in S_{i,t}} e^{\beta_{k,l}}} \right]$$

The above optimization problem is separable in  $\underline{\alpha}$  &  $\underline{\beta}$ . Optimizing over  $\underline{\alpha}$  we get

$$\alpha_k^{(t+1)} = \frac{\sum_i h_{ik}^{(t)}}{\sum_{k'} \sum_i h_{ik'}^{(t)}}$$

$$\beta_k^{(t+1)} = \underset{k}{\operatorname{argmax}} \sum_i h_{ik}^{(t)} \sum_{t \in T_i} \log \frac{e^{\beta_{k,j_i,t}}}{\sum_{l \in S_{i,t}} e^{\beta_{k,l}}} \quad \forall k$$

Implementation note: In order to determine  $\beta_k^{(t+1)}$ , you can use  $\beta_k^{(t)}$  as the initial solution. So, we can do the updates in a "lazy" fashion by running only one MM update step. We can write

$$\beta_{k,j}^{(t+1)} = \beta_{k,j}^{(t)} + \log \frac{\sum_i h_{ik}^{(t)} \sum_{t \in T_i} \mathbb{1}[j = j_{i,t}]}{\sum_i h_{ik}^{(t)} \sum_{t \in T_i} \mathbb{P}_k^{(t)}[j|S_{i,t}]}$$

## Frank-Wolfe algorithm for estimating a rank-based choice model

set up: We have  $n$  items. We have observation of the form

$f_{j,S}$  = fraction of purchases of item  $j$  when  $S$  was offered for a collection of offer sets  $S_1, S_2, \dots, S_m$ .

Model: We assume that the data are generated as follows. The population is described by a generation distribution (PMF). Over all possible rankings/pref lists of the  $n$  items. In particular,  $\lambda_\sigma$  is the probability of sampling  $\sigma$ , where  $\sum_\sigma \lambda_\sigma = 1$ ,  $\lambda_\sigma \geq 0 \quad \forall \sigma$ , when given an offer set  $S$ , the customer samples a preference list  $\sigma$  according to  $\underline{\lambda}$  and chooses the most preferred item from  $S$  according to  $\sigma$ .

Estimation: We estimate the model through MLE.

$$\begin{aligned} l(\lambda) &= \sum_{i=1}^m \sum_{j \in S_i} (\log \mathbb{P}_\lambda(j|S_i) f_{j,S_i}) \\ &= \sum_{i=1}^m \sum_{j \in S_i} (\log (\sum_\sigma \lambda_\sigma \mathbb{1}[\sigma(j) < \sigma(k) \quad \forall k \in S_i \setminus \{j\}]) f_{j,S_i}) \\ &\equiv j \text{ is most preferred among items in } S_i \text{ under } \sigma \\ &= \sum_{i=1}^m \sum_{j \in S_i} ((\log \sum_\sigma \lambda_\sigma \mathbb{1}[\sigma(j) < \sigma(k) \quad \forall k \in S_i \setminus \{j\}]) f_{j,S_i}) \end{aligned}$$

The MLE problem now become

$$\begin{aligned} \max_{\underline{\lambda}} \quad & \sum_{i=1}^m \sum_{j \in S_i} f_{j,S_i} \log(\sum_\sigma \lambda_\sigma \mathbb{1}[\sigma(j) < \sigma(k) \quad \forall k \in S_i \setminus \{j\}]) \\ \text{s.t.} \quad & \sum_\sigma \lambda_\sigma = 1 \quad \lambda_\sigma \geq 0 \quad \forall \sigma \end{aligned}$$

Because the objective is concave in the variable  $\lambda_\sigma$  and the constraints are linear, the above optimization problem is a convex program, albeit a large dimensional one.

Remarks:

1. In general, the above program has multiple optima. For tractability reasons, we choose a solution that has a small support size, where the support of  $\underline{\lambda}$  is defined as  $\{\sigma : \lambda_\sigma > 0\}$ .
2. Consider the following reformulation of the optimization problem:

$$\begin{aligned}
& \max_{\underline{\lambda}, \underline{g}} \sum_{i=1}^m \sum_{j \in S_i} f_{j, S_i} \log g_{j, S_i} \\
& \text{s.t. } g_{j, S_i} = \sum_{\sigma} \lambda_{\sigma} \mathbb{1}[\sigma; j; S_i] \quad \forall j \in S_i, \quad i = 1, \dots, m \\
& \sum_{\sigma} \lambda_{\sigma} = 1, \quad \lambda_{\sigma} \geq 0 \quad \forall \sigma
\end{aligned}$$

Now consider the following vectorization. Let  $L = \sum_{i=1}^m |S_i|$  and  $\underline{g} \in \mathbb{R}^L$  s.t.  $(\underline{g})_{j, S_i} = g_{j, S_i}$ . Also, for any  $\sigma$ , let  $\underline{e}_{\sigma} \in \{0, 1\}^L$  s.t.  $(\underline{e}_{\sigma})_{j, S_i} = \mathbb{1}[\sigma; j; S_i]$ .

We can rewrite the above optimization problem as

$$\begin{aligned}
& \max_{\underline{g}} \sum_{i=1}^m \sum_{j \in S_i} f_{j, S_i} \log g_{j, S_i} \\
& \text{s.t. } \underline{g} \in \text{conv}(\{\underline{e}_{\sigma} : \sigma\})
\end{aligned}$$

Remarks:

Suppose  $\underline{g}^*$  is an optimal solution to the above program. It follows from Caratheodory's Theorem that  $\exists$  a convex decomposition of  $\underline{g}^*$  in terms of  $\underline{e}_{\sigma}$  with support size at most  $L + 1$ .